#### **Computer Architecture** ELE 475 / COS 475 Slide Deck 14: Interconnection Networks David Wentzlaff **Department of Electrical Engineering Princeton University**





#### Overview of Interconnection Networks: Buses



#### Overview of Interconnection Networks: Buses



#### Overview of Interconnection Networks: Buses



#### Overview of Interconnection Networks: Point-to-point / Switched



#### Overview of Interconnection Networks: Point-to-point / Switched



#### Explicit Message Passing (Programming)

- Send(Destination, \*Data)
- Receive(&Data)
- Receive(Source, & Data)

- Unicast (one-to-one)
- Multicast (one-to-multiple)
- Broadcast (one-to-all)

#### Message Passing Interface (MPI)

```
#include <stdio.h>
#include <assert.h>
#include <mpi.h>
int main (int argc, char **argv) {
  int myid, numprocs, x, y;
  int tag = 475;
  MPI Status status;
  MPI Init(&argc,&argv);
  MPI Comm size(MPI COMM WORLD,&numprocs);
  MPI Comm rank(MPI COMM WORLD,&myid);
  assert(numprocs == 2);
  if(myid==0) {
    x = 475;
    MPI Send(&x, 1, MPI INT, 1, tag, MPI COMM WORLD);
    MPI Recv(&y, 1, MPI INT, 1, tag, MPI COMM WORLD, &status);
    printf("received number: ELE %d A\n", y);
  }
  else {
    MPI Recv(&y, 1, MPI INT, 0, tag, MPI COMM WORLD, &status);
    v += 105;
    MPI Send(&y, 1, MPI INT, 0, tag, MPI COMM WORLD);
  }
  MPI Finalize();
  exit(0);
}
```

#### Message Passing vs. Shared Memory

- Message Passing
  - Memory is private
  - Explicit send/receive to communicate
  - Message contains data and synchronization
  - Need to know Destination on generation of data (send)
  - Easy for Producer-Consumer
- Shared Memory
  - Memory is shared
  - Implicit communication via loads and stores
  - Implicit synchronization needed via Fences, Locks, and Flags
  - No need to know Destination on generation of date (can store in memory and user of data can pick up later)
  - Easy for multiple threads accessing a shared table
  - Needs Locks and critical sections to synchronize access

#### Shared Memory Tunneled over Messaging

• Software

Turn loads and stores into sends and receives

- Hardware
  - Replace bus communications with messages sent between cores and between cores and memory



#### Shared Memory Tunneled over Messaging

• Software

Turn loads and stores into sends and receives

- Hardware
  - Replace bus communications with messages sent between cores and between cores and memory



#### Messaging Tunneled over Shared Memory

 Use software queues (FIFOs) with locks to transmit data directly between cores by loads and stores to memory



#### Interconnect Design

- Switching
- Topology
- Routing
- Flow Control



- Flit: flow control digit (Basic unit of flow control)
- Phit: physical transfer digit (Basic unit of data transferred in one clock)

### Switching

- Circuit Switched
- Store and Forward
- Cut-through
- Wormhole

Bus



Pipelined Bus / Segmented Bus



Ring / IP Torus





# 2D Mesh 0 6

Topology

2D

Torus



Star/ Fully connected crossbar



Omega Network







3- gry 3- cube mesh

#### **Topology Parameters**

- Routing Distance: Number of links between two points
- Diameter: Maximum routing distance between any two points
- Average Distance
- Minimum Bisection Bandwidth (Bisection Bandwidth): The bandwidth of a minimal cut though the network such that the network is divided into two sets of nodes
- Degree of a Router

#### **Topology Parameters**

2D Mesh

Diameter:  $2\sqrt{N}$  - 2

Bisection Bandwidth:  $2\sqrt{N}$ 

Degree of a Router: 5



#### **Topology Influenced by Packaging**

Star/ Fully connected crossbar . Wiring grows as N-1



 Physically hard to pack into 3-space (pack in sphere?)

#### Topology Influenced by Packaging

- Packing N dimensions in N-1 space leads to long wires
- Packing N dimensions in N-2 space leads to really long wires



#### Network Performance

- Bandwidth: The rate of data that can be transmitted over the network (network link) in a given time
- Latency: The time taken for a message to be sent from sender to receiver
- Bandwidth can affect latency
  - Reduce congestion
  - Messages take fewer Flits and Phits
- Latency can affect Bandwidth
  - Round trip communication can be limited by latency
  - Round trip flow-control can be limited by latency

Sprializer deserializer - Router Pipeline: RORIR2 -Link traversal: LO LI Packet I Head Phit SRORIRZ LOLIRORIR2D Body Phit SRORIRZ LOLIRORIR2D Body Phit S RO RI AZ LOLI RO RI AZ D Toil Phit S RORI R2 LOLI RORI R2D Sevialization Channe Ratency Router Pipeline Latency Latency to

#### Anatomy of Message Latency

 $T = T_{head} + L/b$ 

 $T_{head}$ : Head Phit Latency, includes  $t_{C}$ ,  $t_{R}$ , hop count, and contention

Unloaded Latency:  $T_0 = H_R * t_R + H_C * t_C + L/b$ 

#### Anatomy of Message Latency

Packet 1 Head Phit S RO RIR2 LOLI RO RIR2D Body Phit S RO RI R2 LOLI RO RI R2D Body Phit S RO RI R2 LOLI RO RI R2D Toil Phit S RO RI R2 LOLI RO RI R2D Serialization Channel Latency Router Pipeline Latency to Latency to

**Unloaded Latency:** 

 $T_0 = H_R * t_R + H_C * t_C + L/b$ 

Shorter routes

Faster channels Faster routers

Wider channels or shorter messages

#### **Interconnection Network Performance**



### Routing

- Oblivious (routing path independent of state of network)
  - Deterministic
  - Non-Deterministic
- Adaptive (routing path depends on state of network)

#### **Flow Control**

- Local (Link or hop based) Flow Control
- End-to-end (Long distance)

#### Deadlock

 Deadlock can occur if cycle possible in "Waitsfor" graph



## Deadlock Example (Waits-for and Holds analysis)





#### Deadlock Avoidance vs. Deadlock Recovery

Deadlock Avoidance

Protocol designed to never deadlock

- Deadlock Recovery
  - Allow Deadlock to occur and then resolve deadlock usually through use of more buffering

#### Acknowledgements

- These slides contain material developed and copyright by:
  - Arvind (MIT)
  - Krste Asanovic (MIT/UCB)
  - Joel Emer (Intel/MIT)
  - James Hoe (CMU)
  - John Kubiatowicz (UCB)
  - David Patterson (UCB)
  - Christopher Batten (Cornell)
- MIT material derived from course 6.823
- UCB material derived from course CS252 & CS152
- Cornell material derived from course ECE 4750

On/off with Combinational Stall Signal



stall D Packet A ₿ CD unstall D ABCD 2 ABC(D)D D 3 ABCCCCC ABBBBBCD AAAAB CD BC D В



B A CD A B B A 5 CD BB B AAAAA B B C D .



ABCD ABCD ABC DDDDD ABCQ ABCQ ABBBBBC D AAAABCD BCD

## Credit - Based Flow Control



- LA CREDIT COUNTER
  - · Decrement counter on send packet
  - · Increment counter on credit recieved

#### Copyright © 2013 David Wentzlaff